

```

1
2 from a0_items import *
3
4 # pip install requests-html
5 from requests_html import HTMLSession
6
7
8
9 def get_proxy_list ():
10
11     LIST_proxy_ip_port = []
12
13     #get GOOD proxy from DB_proxies
14
15
16     return LIST_proxy_ip_port
17
18
19
20 #with random proxy
21 def get_url_status_code (url):
22
23     LIST_proxy_ip_port = get_proxy_list()
24     random_proxy = random.choice(LIST_proxy_ip_port + [''])
25
26     x = requests.get(url, proxies=random_proxy)
27
28     status_code = x.status_code
29     #print(x.status_code)
30
31     return status_code
32
33
34 #print(get_url_status_code (url='https://w3schools.com'))
35
36
37
38 #https://www.scrapingbee.com/blog/python-requests-proxy/
39 def request_url_with_session (url):
40
41     session = requests.Session()
42
43     LIST_proxy_ip_port = get_proxy_list()
44     random_proxy = random.choice(LIST_proxy_ip_port + [''])
45
46     session.proxies = random_proxy
47
48     response = session.get(url)
49     print (response)
50
51
52 request_url_with_session (url='https://w3schools.com')
53
54
55
56
57 #
58 https://stackoverflow.com/questions/26393231/using-python-requests-with-javascript-pages
59 def get_page_with_JS ():
60
61     session = HTMLSession()
62
63     r = session.get('https://duckduckgo.com/?va=k&t=ht&q=top+RSS+feeds+china&ia=web')
64
65     page = r.html.render() # this call executes the js in the page
66
67     print (page)

```

```

67
68 #get_page_with_JS ()
69
70
71
72 def get_ALL_urls (url):
73
74     try:
75
76         reqs = requests.get(url)
77
78         soup = BeautifulSoup(reqs.text, 'html.parser')
79
80         LIST_url = []
81
82         for link in soup.findAll('a'):
83
84             #print(link.get('href'))
85             LIST_url.append(link.get('href'))
86
87         return LIST_url
88
89     except requests.exceptions.RequestException as errex:
90         print("Exception request")
91
92
93
94
95 '''
96 LIST_url = get_ALL_urls (url= 'https://rss.feedspot.com/algeria_rss_feeds/')
97
98 for x in LIST_url:
99
100     if 'rss' in x or 'feed' in x and 'http' in x and 'feedspot' not in x:
101
102         print (x)
103
104     else:
105
106         print ('rss NOT in')
107 '''
108
109
110 def get_Image_urls (url):
111
112     LIST_image_url = []
113
114     response = requests.get(url)
115     soup = BeautifulSoup(response.content, "html.parser")
116
117     for img in soup.findAll('img'):
118
119         src = img.get("src")
120
121         if src:
122
123             # resolve any relative urls to absolute urls using base URL
124             src = requests.compat.urljoin(url, src)
125             #print(src)
126
127             LIST_image_url.append (src)
128
129     return LIST_image_url
130
131
132 #print (get_Image_urls ('https://en.wikipedia.org/wiki/Main_Page'))
133

```

```
134
135
136
137 def create_url_with_params (url, params):
138
139     r = requests.get(url, params)
140     url_with_params = r.url
141     print(url_with_params)
142
143     return url_with_params
144
145
146 # should only have TFM_Task, TFM_task_action parameters,
147 # Other values should use TFMdata1, 2..3 as value keys, so that PHP only needs to GET
148 # TFM_task & TFM_task variables
149 parameters = {'TFM_task': 'MSQM_Country_language',
150               'TFM_task_action':
151                 'Folders_CREATE_TFMdata1_/data/raw/_TFMdata2_xyzfile.db_TFMdata2_helloFILE.
152                 txt',
153               }
154
155 #print (create_url_with_params (url='https://httpbin.org/get', params=parameters))
```